



## Longest Common Extensions in Sublinear Space

**Bille, Philip; Gørtz, Inge Li; Knudsen, Mathias Bæk Tejs; Lewenstein, Moshe; Vildhøj, Hjalte Wedel**

*Published in:*  
Combinatorial Pattern Matching

*Link to article, DOI:*  
[10.1007/978-3-319-19929-0\\_6](https://doi.org/10.1007/978-3-319-19929-0_6)

*Publication date:*  
2015

*Document Version*  
Peer reviewed version

[Link back to DTU Orbit](#)

*Citation (APA):*  
Bille, P., Gørtz, I. L., Knudsen, M. B. T., Lewenstein, M., & Vildhøj, H. W. (2015). Longest Common Extensions in Sublinear Space. In *Combinatorial Pattern Matching: 26th Annual Symposium, CPM 2015, Ischia Island, Italy, June 29 -- July 1, 2015, Proceedings* (pp. 65-76). Springer. Lecture Notes in Computer Science Vol. 9133  
[https://doi.org/10.1007/978-3-319-19929-0\\_6](https://doi.org/10.1007/978-3-319-19929-0_6)

---

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

# Longest Common Extensions in Sublinear Space

Philip Bille<sup>1 \*</sup>, Inge Li Gørtz<sup>1 \*</sup>, Mathias Bæk Tejs Knudsen<sup>2 \*\*</sup>,  
Moshe Lewenstein<sup>3 \*\*\*</sup>, and Hjalte Wedel Vildhøj<sup>1</sup>

<sup>1</sup> Technical University of Denmark, DTU Compute

<sup>2</sup> Department of Computer Science, University of Copenhagen

<sup>3</sup> Bar Ilan University

**Abstract.** The *longest common extension problem* (LCE problem) is to construct a data structure for an input string  $T$  of length  $n$  that supports  $\text{LCE}(i, j)$  queries. Such a query returns the length of the longest common prefix of the suffixes starting at positions  $i$  and  $j$  in  $T$ . This classic problem has a well-known solution that uses  $\mathcal{O}(n)$  space and  $\mathcal{O}(1)$  query time. In this paper we show that for any trade-off parameter  $1 \leq \tau \leq n$ , the problem can be solved in  $\mathcal{O}(\frac{n}{\tau})$  space and  $\mathcal{O}(\tau)$  query time. This significantly improves the previously best known time-space trade-offs, and almost matches the best known time-space product lower bound.

## 1 Introduction

Given a string  $T$ , the *longest common extension* of suffix  $i$  and  $j$ , denoted  $\text{LCE}(i, j)$ , is the length of the longest common prefix of the suffixes of  $T$  starting at position  $i$  and  $j$ . The *longest common extension problem* (LCE problem) is to preprocess  $T$  into a compact data structure supporting fast longest common extension queries.

The LCE problem is a basic primitive that appears as a central subproblem in a wide range of string matching problems such as approximate string matching and its variations [1, 4, 11, 13, 16], computing exact or approximate repetitions [6, 12, 14], and computing palindromes [10, 15]. In many cases the LCE problem is the computational bottleneck.

Here we study the time-space trade-offs for the LCE problem, that is, the space used by the preprocessed data structure vs. the worst-case time used by LCE queries. The input string is given in read-only memory and is not counted in the space complexity. Throughout the paper we use  $\ell$  as a shorthand for  $\text{LCE}(i, j)$ . The standard trade-offs are as follows: At one extreme we can store

---

\* Supported by the Danish Research Council and the Danish Research Council under the Sapere Aude Program (DFF 4005-00267)

\*\* Research partly supported by Mikkel Thorup's Advanced Grant from the Danish Council for Independent Research under the Sapere Aude research career programme and the FNU project AlgoDisc - Discrete Mathematics, Algorithms, and Data Structures.

\*\*\* This research was supported by a Grant from the GIF, the German-Israeli Foundation for Scientific Research and Development, and by a BSF grant 2010437

a suffix tree combined with an efficient nearest common ancestor (NCA) data structure [7,17]. This solution uses  $\mathcal{O}(n)$  space and supports LCE queries in  $\mathcal{O}(1)$  time. At the other extreme we do not store any data structure and instead answer queries simply by comparing characters from left-to-right in  $T$ . This solution uses  $\mathcal{O}(1)$  space and answers an  $\text{LCE}(i, j)$  query in  $\mathcal{O}(\ell) = \mathcal{O}(n)$  time. Recently, Bille et al. [2] presented a number of results. For a trade-off parameter  $\tau$ , they gave: 1) a deterministic solution with  $\mathcal{O}(\frac{n}{\tau})$  space and  $\mathcal{O}(\tau^2)$  query time, 2) a randomized Monte Carlo solution with  $\mathcal{O}(\frac{n}{\tau})$  space and  $\mathcal{O}(\tau \log(\frac{\ell}{\tau})) = \mathcal{O}(\tau \log(\frac{n}{\tau}))$  query time, where all queries are correct with high probability, and 3) a randomized Las Vegas solution with the same bounds as 2) but where all queries are guaranteed to be correct. Bille et al. [2] also gave a lower bound showing that any data structure for the LCE problem must have a time-space product of  $\Omega(n)$  bits.

*Our Results* Let  $\tau$  be a trade-off parameter. We present four new solutions with the following improved bounds. Unless otherwise noted the space bound is the number of words on a standard RAM with logarithmic word size, not including the input string, which is given in read-only memory.

- A deterministic solution with  $\mathcal{O}(n/\tau)$  space and  $\mathcal{O}(\tau \log^2(n/\tau))$  query time.
- A randomized Monte Carlo solution with  $\mathcal{O}(n/\tau)$  space and  $\mathcal{O}(\tau)$  query time, such that all queries are correct with high probability.
- A randomized Las Vegas solution with  $\mathcal{O}(n/\tau)$  space and  $\mathcal{O}(\tau)$  query time.
- A derandomized version of the Monte Carlo solution with  $\mathcal{O}(n/\tau)$  space and  $\mathcal{O}(\tau)$  query time.

Hence, we obtain the first trade-off for the LCE problem with a linear time-space product in the full range from constant to linear space. This almost matches the

	Data Structure			Preprocessing		Trade-off range	Reference
	Space	Query	Correct	Space	Time		
Deterministic	1	$\ell$	always	1	1	-	Store nothing
	$n$	1	always	$n$	$n$	-	Suffix tree + NCA
	$\frac{n}{\tau}$	$\tau^2$	always	$\frac{n}{\tau}$	$\frac{n^2}{\tau}$	$1 \leq \tau \leq \sqrt{n}$	[2]
	$\frac{n}{\tau}$	$\tau \log^2 \frac{n}{\tau}$	always	$\frac{n}{\tau}$	$n^2$	$\frac{1}{\log n} \leq \tau \leq n$	this paper, Sec. 2
	$\frac{n}{\tau}$	$\tau$	always	$\frac{n}{\tau}$	$n^{2+\varepsilon}$	$1 \leq \tau \leq n$	this paper, Sec. 4
Monte Carlo	$\frac{n}{\tau}$	$\tau \log \frac{\ell}{\tau}$	w.h.p.	$\frac{n}{\tau}$	$n$	$1 \leq \tau \leq n$	[2]
	$\frac{n}{\tau}$	$\tau$	w.h.p.	$\frac{n}{\tau}$	$n \log \frac{n}{\tau}$	$1 \leq \tau \leq n$	this paper, Sec. 3
Las Vegas	$\frac{n}{\tau}$	$\tau \log \frac{\ell}{\tau}$	always	$\frac{n}{\tau}$	$n(\tau + \log n)$ w.h.p.	$1 \leq \tau \leq n$	[2]
	$\frac{n}{\tau}$	$\tau$	always	$\frac{n}{\tau}$	$n^{3/2}$ w.h.p.	$1 \leq \tau \leq n$	this paper, Sec. 3.5

**Table 1.** Overview of solutions for the LCE problem. Here  $\ell = \text{LCE}(i, j)$ ,  $\varepsilon > 0$  is an arbitrarily small constant and w.h.p. (with high probability) means with probability at least  $1 - n^{-c}$  for an arbitrarily large constant  $c$ . The data structure is *correct* if it answers all LCE queries correctly.

time-space product lower bound of  $\Omega(n)$  bits, and improves the best deterministic upper bound by a factor of  $\tau$ , and the best randomized bound by a factor  $\log(\frac{n}{\tau})$ . See the columns marked *Data Structure* in Table 1 for a complete overview.

While our main focus is the space and query time complexity, we also provide efficient *preprocessing* algorithms for building the data structures, supporting independent trade-offs between the preprocessing time and preprocessing space. See the columns marked *Preprocessing* in Table 1.

To achieve our results we develop several new techniques and specialized data structures which are likely of independent interest. For instance, in our deterministic solution we develop a novel recursive decomposition of LCE queries and for the randomized solution we develop a new sampling technique for Karp-Rabin fingerprints that allow fast LCE queries. We also give a general technique for efficiently derandomizing algorithms that rely on “few” or “short” Karp-Rabin fingerprints, and apply the technique to derandomize our Monte Carlo algorithm. To the best of our knowledge, this is the first derandomization technique for Karp-Rabin fingerprints.

*Preliminaries* We assume an integer alphabet, i.e.,  $T$  is chosen from some alphabet  $\Sigma = \{0, \dots, n^c\}$  for some constant  $c$ , so every character of  $T$  fits in  $\mathcal{O}(1)$  words. For integers  $a \leq b$ ,  $[a, b]$  denotes the range  $\{a, a+1, \dots, b\}$  and we define  $[n] = [0, n-1]$ . For a string  $S = S[1]S[2] \dots S[|S|]$  and positions  $1 \leq i \leq j \leq |S|$ ,  $S[i \dots j] = S[i]S[i+1] \dots S[j]$  is a *substring* of length  $j-i+1$ ,  $S[i \dots] = S[i, |S|]$  is the  $i^{\text{th}}$  *suffix* of  $S$ , and  $S[\dots i] = S[1, i]$  is the  $i^{\text{th}}$  *prefix* of  $S$ .

## 2 Deterministic Trade-Off

Here we describe a completely deterministic trade-off for the LCE problem with  $\mathcal{O}(\frac{n}{\tau} \log \frac{n}{\tau})$  space and  $\mathcal{O}(\tau \log \frac{n}{\tau})$  query time for any  $\tau \in [1, n]$ . Substituting  $\hat{\tau} = \tau / \log(n/\tau)$ , we obtain the bounds reflected in Table 1 for  $\hat{\tau} \in [1/\log n, n]$ .

A key component in this solution is the following observation that allows us to reduce an  $\text{LCE}(i, j)$  query on  $T$  to another query  $\text{LCE}(i', j')$  where  $i'$  and  $j'$  are both indices in either the first or second half of  $T$ .

**Observation 1** *Let  $i, j$  and  $j'$  be indices of  $T$ , and suppose that  $\text{LCE}(j', j) \geq \text{LCE}(i, j)$ . Then  $\text{LCE}(i, j) = \min(\text{LCE}(i, j'), \text{LCE}(j', j))$ .*

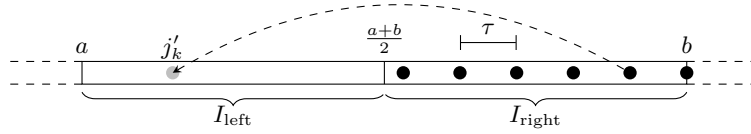
We apply Observation 1 recursively to bring the indices of the initial query within distance  $\tau$  in  $\mathcal{O}(\log(n/\tau))$  rounds. We show how to implement each round with a data structure using  $\mathcal{O}(n/\tau)$  space and  $\mathcal{O}(\tau)$  time. This leads to a solution using  $\mathcal{O}(\frac{n}{\tau} \log \frac{n}{\tau})$  space and  $\mathcal{O}(\tau \log \frac{n}{\tau})$  query time. Finally in Section 2.4, we show how to efficiently solve the LCE problem for indices within distance  $\tau$  in  $\mathcal{O}(n/\tau)$  space and  $\mathcal{O}(\tau)$  time by exploiting periodicity properties of LCEs.

### 2.1 The Data Structure

We will store several data structures, each responsible for a specific subinterval  $I = [a, b] \subseteq [1, n]$  of positions of the input string  $T$ . Let  $I_{\text{left}} = [a, (a+b)/2]$ ,

$I_{\text{right}} = ((a+b)/2, b]$ , and  $|I| = b - a + 1$ . The task of the data structure for  $I$  will be to reduce an  $\text{LCE}(i, j)$  query where  $i, j \in I$  to one where both indices belong to either  $I_{\text{left}}$  or  $I_{\text{right}}$ .

The data structure stores information for  $\mathcal{O}(|I|/\tau)$  suffixes of  $T$  that start in  $I_{\text{right}}$ . More specifically, we store information for the suffixes starting at positions  $b - k\tau \in I_{\text{right}}$ ,  $k = 0, 1, \dots, (|I|/2)/\tau$ . We call these the *sampled positions* of  $I_{\text{right}}$ . See Figure 1 for an illustration.



**Fig. 1.** Illustration of the contents of the data structure for the interval  $I = [a, b]$ . The black dots are the sampled positions in  $I_{\text{right}}$ , and each such position has a pointer to an index  $j'_k \in I_{\text{left}}$ .

For every sampled position  $b - k\tau \in I_{\text{right}}$ ,  $k = 0, 1, \dots, (|I|/2)/\tau$ , we store the index  $j'_k$  of the suffix starting in  $I_{\text{left}}$  that achieves the maximum LCE value with the suffix starting at the sampled position, i.e.,  $T[b - k\tau\dots]$  (ties broken arbitrarily). Along with  $j'_k$ , we also store the value of the LCE between suffix  $T[j'_k\dots]$  and  $T[b - k\tau\dots]$ . Formally,  $j'_k$  and  $L_k$  are defined as follows,

$$j'_k = \operatorname{argmax}_{h \in I_{\text{left}}} \text{LCE}(h, b - k\tau) \quad \text{and} \quad L_k = \text{LCE}(j'_k, b - k\tau).$$

**Building the structure** We construct the above data structure for the interval  $[1, n]$ , and build it recursively for  $[1, n/2]$  and  $(n/2, n]$ , stopping when the length of the interval becomes smaller than  $\tau$ .

## 2.2 Answering a Query

We now describe how to reduce a query  $\text{LCE}(i, j)$  where  $i, j \in I$  to one where both indices are in either  $I_{\text{left}}$  or  $I_{\text{right}}$ . Suppose without loss of generality that  $i \in I_{\text{left}}$  and  $j \in I_{\text{right}}$ . We start by comparing  $\delta < \tau$  pairs of characters of  $T$ , starting with  $T[i] = T[j]$ , until 1) we encounter a mismatch, 2) both positions are in  $I_{\text{right}}$  or 3) we reach a sampled position in  $I_{\text{right}}$ . It suffices to describe the last case, in which  $T[i, i + \delta] = T[j, j + \delta]$ ,  $i + \delta \in I_{\text{left}}$  and  $j + \delta = b - k\tau \in I_{\text{right}}$  for some  $k$ . Then by Observation 1, we have that

$$\begin{aligned} \text{LCE}(i, j) &= \delta + \text{LCE}(i + \delta, j + \delta) \\ &= \delta + \min(\text{LCE}(i + \delta, j'_k), \text{LCE}(j'_k, b - k\tau)) \\ &= \delta + \min(\text{LCE}(i + \delta, j'_k), L_k). \end{aligned}$$

Thus, we have reduced the original query to computing the query  $\text{LCE}(i + \delta, j'_k)$  in which both indices are in  $I_{\text{left}}$ .

### 2.3 Analysis

Each round takes  $\mathcal{O}(\tau)$  time and halves the upper bound for  $|i - j|$ , which initially is  $n$ . Thus, after  $\mathcal{O}(\tau \log(n/\tau))$  time, the initial LCE query has been reduced to one where  $|i - j| \leq \tau$ . At each of the  $\mathcal{O}(\log(n/\tau))$  levels, the number of sampled positions is  $(n/2)/\tau$ , so the total space used is  $\mathcal{O}((n/\tau) \log(n/\tau))$ .

### 2.4 Queries with Nearby Indices

We now describe the data structure used to answer a query  $\text{LCE}(i, j)$  when  $|i - j| \leq \tau$ . We first give some necessary definitions and properties of periodic strings. We say that the integer  $1 \leq p \leq |S|$  is a *period* of a string  $S$  if any two characters that are  $p$  positions apart in  $S$  match, i.e.,  $S[i] = S[i + p]$  for all positions  $i$  s.t.  $1 \leq i < i + p \leq |S|$ . The following is a well-known property of periods.

**Lemma 1 (Fine and Wilf [5]).** *If a string  $S$  has periods  $a$  and  $b$  and  $|S| \geq |a| + |b| - \gcd(a, b)$ , then  $\gcd(a, b)$  is also a period of  $S$ .*

The *period* of  $S$  is the smallest period of  $S$  and we denote it by  $\text{per}(S)$ . If  $\text{per}(S) \leq |S|/2$ , we say  $S$  is *periodic*. A periodic string  $S$  might have many periods smaller than  $|S|/2$ , however it follows from the above lemma that

**Corollary 1.** *All periods smaller than  $|S|/2$  are multiples of  $\text{per}(S)$ .*

**The Data Structure** Let  $T_k = T[k\tau \dots (k+2)\tau - 1]$  denote the substring of length  $2\tau$  starting at position  $k\tau$  in  $T$ ,  $k = 0, 1, \dots, n/\tau$ . For the strings  $T_k$  that are periodic, let  $p_k = \text{per}(T_k)$  be the period. For every periodic  $T_k$ , the data structure stores the length  $\ell_k$  of the maximum substring starting at position  $k\tau$ , which has period  $p_k$ . Nothing is stored if  $T_k$  is aperiodic.

**Answering a Query** We may assume without loss of generality that  $i = k\tau$ , for some integer  $k$ . If not, then we check whether  $T[i + \delta] = T[j + \delta]$  until  $i + \delta = k\tau$ . Hence, assume that  $i = k\tau$  and  $j = i + d$  for some  $0 < d \leq \tau$ . In  $\mathcal{O}(\tau)$  time, we first check whether  $T[i + \delta] = T[j + \delta]$  for all  $\delta \in [0, 2\tau]$ . If we find a mismatch we are done, and otherwise we return  $\text{LCE}(i, j) = \ell_k - d$ .

**Correctness** If a mismatch is found when checking that  $T[i + \delta] = T[j + \delta]$  for all  $\delta \in [0, 2\tau]$ , the answer is clearly correct. Otherwise, we have established that  $d \leq \tau$  is a period of  $T_k$ , so  $T_k$  is periodic and  $d$  is a multiple of  $p_k$  (by Corollary 1). Consequently,  $T[i + \delta] = T[j + \delta]$  for all  $\delta$  s.t.  $d + \delta \leq \ell_k$ , and thus  $\text{LCE}(i, j) = \ell_k - d$ .

## 2.5 Preprocessing

The preprocessing details appear in the full version of this paper [3].

## 3 Randomized Trade-Offs

In this section we describe a randomized LCE data structure using  $\mathcal{O}(n/\tau)$  space with  $\mathcal{O}(\tau + \log \frac{\ell}{\tau})$  query time. In Section 3.6 we describe another  $\mathcal{O}(n/\tau)$ -space LCE data structure that either answers an LCE query in constant time, or provides a certificate that  $\ell \leq \tau^2$ . Combining the two data structures, shows that the LCE problem can be solved in  $\mathcal{O}(n/\tau)$  space and  $\mathcal{O}(\tau)$  time.

The randomization comes from our use of Karp-Rabin fingerprints [9] for comparing substrings of  $T$  for equality. Before describing the data structure, we start by briefly recapping the most important definitions and properties of Karp-Rabin fingerprints.

### 3.1 Karp-Rabin Fingerprints

For a prime  $p$  and  $x \in [p]$  the Karp-Rabin fingerprint [9], denoted  $\phi_{p,x}(T[i\dots j])$ , of the substring  $T[i\dots j]$  is defined as

$$\phi_{p,x}(T[i\dots j]) = \sum_{i \leq k \leq j} T[k]x^{k-i} \bmod p .$$

If  $T[i\dots j] = T[i'\dots j']$  then clearly  $\phi_{p,x}(T[i\dots j]) = \phi_{p,x}(T[i'\dots j'])$ . In the Monte Carlo and the Las Vegas algorithms we present we will choose  $p$  such that  $p = \Theta(n^{4+c})$  for some constant  $c > 0$  and  $x$  uniformly from  $[p] \setminus \{0\}$ . In this case a simple union bound shows that the converse is also true with high probability, i.e.,  $\phi$  is *collision-free* on all substring pairs of  $T$  with probability at least  $1 - n^{-c}$ . Storing a fingerprint requires  $\mathcal{O}(1)$  space. When  $p, x$  are clear from the context we write  $\phi = \phi_{p,x}$ .

For shorthand we write  $f(i) = \phi(T[1, i])$ ,  $i \in [1, n]$  for the fingerprint of the  $i^{\text{th}}$  prefix of  $T$ . Assuming that we store the exponent  $x^i \bmod p$  along with the fingerprint  $f(i)$ , the following two properties of fingerprints are well-known and easy to show.

**Lemma 2.** 1) Given  $f(i)$ , the fingerprint  $f(i \pm a)$  for some integer  $a$ , can be computed in  $\mathcal{O}(a)$  time. 2) Given fingerprints  $f(i)$  and  $f(j)$ , the fingerprint  $\phi(T[i\dots j])$  can be computed in  $\mathcal{O}(1)$  time.

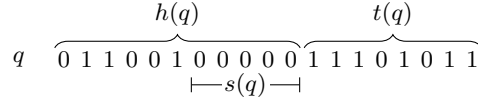
In particular this implies that for a fixed length  $l$ , the fingerprint of all substrings of length  $l$  of  $T$  can be enumerated in  $\mathcal{O}(n)$  time using a sliding window.

### 3.2 Overview

The main idea in our solution is to binary search for the  $\text{LCE}(i, j)$  value using Karp-Rabin fingerprints. Suppose for instance that  $\phi(T[i, i + M]) \neq \phi(T[j, j + M])$  for some integer  $M$ , then we know that  $\text{LCE}(i, j) \leq M$ , and thus we can find the true  $\text{LCE}(i, j)$  value by comparing  $\log(M)$  additional pair of fingerprints. The challenge is to obtain the fingerprints quickly when we are only allowed to use  $\mathcal{O}(n/\tau)$  space. We will partition the input string  $T$  into  $n/\tau$  blocks each of length  $\tau$ . Within each block we sample a number of equally spaced positions. The data structure consists of the fingerprints of the prefixes of  $T$  that ends at the sampled positions, i.e., we store  $f(i)$  for all sampled positions  $i$ . In total we sample  $\mathcal{O}(n/\tau)$  positions. If we just sampled a single position in each block (similar to the approach in [2]), we could compute the fingerprint of any substring in  $\mathcal{O}(\tau)$  time (see Lemma 2), and the above binary search algorithm would take time  $\mathcal{O}(\tau \log n)$  time. We present a new sampling technique that only samples an additional  $\mathcal{O}(n/\tau)$  positions, while improving the query time to  $\mathcal{O}(\tau + \log(\ell/\tau))$ .

**Preliminary definitions** We partition the input string  $T$  into  $n/\tau$  blocks of  $\tau$  positions, and by *block*  $k$  we refer to the positions  $[k\tau, k\tau + \tau)$ , for  $k \in [n/\tau]$ .

We assume without loss of generality that  $n$  and  $\tau$  are both powers of two. Every position  $q \in [1, n]$  can be represented as a bit string of length  $\lg n$ . Let  $q \in [1, n]$  and consider the binary representation of  $q$ . We define the leftmost  $\lg(n/\tau)$  bits and rightmost  $\lg(\tau)$  bits to be the *head*, denoted  $h(q)$  and the *tail*, denoted  $t(q)$ , respectively. A position is *block aligned* if  $t(q) = 0$ . The *significance* of  $q$ , denoted  $s(q)$ , is the number of trailing zeros in  $h(q)$ . Note that the  $\tau$  positions in any fixed block  $k \in [n/\tau]$  all have the same head, and thus also the same significance, which we denote by  $\mu_k$ . See Figure 2.



**Fig. 2.** Example of the definitions for the position  $q = 205035$  in a string of length  $n = 2^{19}$  with block length  $\tau = 2^8$ . Here  $h(q)$  is the first  $\lg(n/\tau) = 11$  bits, and  $t(q)$  is the last  $\lg(\tau) = 8$  bits in the binary representation of  $q$ . The significance is  $s(q) = 5$ .

### 3.3 The Monte Carlo Data Structure


The data structure consists of the values  $f(i)$ ,  $i \in \mathcal{S}$ , for a specific set of sampled positions  $\mathcal{S} \subseteq [1, n]$ , along with the information necessary in order to look up the values in constant time. We now explain how to construct the set  $\mathcal{S}$ . In block  $k \in [n/\tau]$  we will sample  $b_k = \min \{2^{\lfloor \mu_k/2 \rfloor}, \tau\}$  evenly spaced positions, where  $\mu_k$  is the significance of the positions in block  $k$ , i.e.,  $\mu_k = s(k\tau)$ . More



precisely, in block  $k$  we sample the positions  $\mathcal{B}_k = \{k\tau + j\tau/b_k \mid j \in [b_k]\}$ , and let  $\mathcal{S} = \cup_{k \in [n/\tau]} \mathcal{B}_k$ . See Figure 3.

We now bound the size of  $\mathcal{S}$ . The significance of a block is at most  $\lg(n/\tau)$ , and there are exactly  $2^{\lg(n/\tau) - \mu}$  blocks with significance  $\mu$ , so

$$|\mathcal{S}| = \sum_{k=0}^{n/\tau-1} b_k \leq \sum_{\mu=0}^{\lg(n/\tau)} 2^{\lg(n/\tau) - \mu} 2^{\lfloor \mu/2 \rfloor} \leq \frac{n}{\tau} \sum_{\mu=0}^{\infty} 2^{-\mu/2} = (2 + \sqrt{2}) \frac{n}{\tau} = \mathcal{O}\left(\frac{n}{\tau}\right).$$

$k$	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$\mu_k$	4	0	1	0	2	0	1	0	3	0	1	0	2	0	1	0
$b_k$	4	1	1	1	2	1	1	1	2	1	1	1	2	1	1	1
$T$																

**Fig. 3.** Illustration of a string  $T$  partitioned into 16 blocks each of length  $\tau$ . The significance  $\mu_k$  for the positions in each block  $k \in [n/\tau]$  is shown, as well as the  $b_k$  values. The block dots are the sampled positions  $\mathcal{S}$ .

### 3.4 Answering a query

We now describe how to answer an  $\text{LCE}(i, j)$  query. We will assume that  $i$  is block aligned, i.e.,  $i = k\tau$  for some  $k \in [n/\tau]$ . Note that we can always obtain this situation in  $\mathcal{O}(\tau)$  time by initially comparing at most  $\tau - 1$  pairs of characters of the input string directly.

Algorithm 1 shows the query algorithm. It performs an exponential search to locate the block in which the first mismatch occurs, after which it scans the block directly to locate the mismatch. The search is performed by calls to `check`( $i, j, c$ ), which computes and compares  $\phi(T[i \dots i + c])$  and  $\phi(T[j \dots j + c])$ . In other words, assuming that  $\phi$  is collision-free, `check`( $i, j, c$ ) returns `true` if  $\text{LCE}(i, j) \geq c$  and `false` otherwise.

**Analysis** We now prove that Algorithm 1 correctly computes  $\ell = \text{LCE}(i, j)$  in  $\mathcal{O}(\tau + \log(\ell/\tau))$  time. The algorithm is correct assuming that `check`( $i, j, 2^\mu \tau$ ) always returns the correct answer, which will be the case if  $\phi$  is collision-free.

The following is the key lemma we need to bound the time complexity.

**Lemma 3.** *Throughout Algorithm 1 it holds that  $\ell \geq (2^\mu - 1)\tau$ ,  $s(j) \geq \mu$ , and  $\mu$  is increased in at least every second iteration of the first **while**-loop.*

*Proof.* We first prove that  $s(j) \geq \mu$ . The claim holds initially. In the first loop  $j$  is changed to  $j + 2^\mu \tau$ , and  $s(j + 2^\mu \tau) \geq \min\{s(j), s(2^\mu \tau)\} = \min\{s(j), \mu\} = \mu$ , where the last equality follows from the induction hypothesis  $s(j) \geq \mu$ . Moreover,  $\mu$  is only incremented when  $s(j) > \mu$ . In the second loop  $j$  is changed to  $j + 2^{\mu-1} \tau$ ,

---

**Algorithm 1** Computing the answer to a query  $\text{LCE}(i, j)$ 

---

```
1: procedure  $\text{LCE}(i, j)$ 
2:    $\hat{\ell} \leftarrow 0$ 
3:    $\mu \leftarrow 0$ 
4:   while  $\text{check}(i, j, 2^\mu \tau)$  do  $\triangleright$  Compute an interval such that  $\ell \in [\hat{\ell}, 2\hat{\ell}]$ .
5:      $(i, j, \hat{\ell}) \leftarrow (i + 2^\mu \tau, j + 2^\mu \tau, \hat{\ell} + 2^\mu \tau)$ 
6:     if  $s(j) > \mu$  then
7:        $\mu \leftarrow \mu + 1$ 
8:   while  $\mu > 0$  do  $\triangleright$  Identify the block in which the first mismatch occurs
9:     if  $\text{check}(i, j, 2^{\mu-1} \tau)$  then
10:       $(i, j, \hat{\ell}) \leftarrow (i + 2^{\mu-1} \tau, j + 2^{\mu-1} \tau, \hat{\ell} + 2^{\mu-1} \tau)$ 
11:       $\mu \leftarrow \mu - 1$ 
12:   while  $T[i] = T[j]$  do  $\triangleright$  Scan the final block left to right to find the mismatch
13:      $(i, j, \hat{\ell}) \leftarrow (i + 1, j + 1, \hat{\ell} + 1)$ 
14:   return  $\hat{\ell}$ 
```

---

which under the assumption that  $s(j) \geq \mu$ , has significance  $s(j + 2^{\mu-1} \tau) = \mu - 1$ . Hence the invariant is restored when  $\mu$  is decremented at line 11.

Now consider an iteration of the first loop where  $\mu$  is not incremented, i.e.,  $s(j) = \mu$ . Then  $\frac{j}{2^\mu \tau}$  is an odd integer, i.e.  $\frac{j+2^\mu \tau}{2^\mu \tau}$  is even, and hence  $s(j + 2^\mu \tau) > \mu$ , so  $\mu$  will be incremented in the next iteration of the loop.

In order to prove that  $\ell \geq (2^\mu - 1) \tau$  we will prove that  $\hat{\ell} \geq (2^\mu - 1) \tau$  in the first loop. This is trivial by induction using the observation that  $(2^\mu - 1) \tau + 2^\mu \tau = (2^{\mu+1} - 1) \tau$ .  $\square$

Since  $\ell \geq (2^\mu - 1) \tau$  and  $\mu$  is increased at least in every second iteration of the first loop and decreased in every iteration of the second loop, it follows that there are  $\mathcal{O}(\log(\ell/\tau))$  iterations of the two first loops. The last loop takes  $\mathcal{O}(\tau)$  time. It remains to prove that the time to evaluate the  $\mathcal{O}(\log(\ell/\tau))$  calls to  $\text{check}(i, j, 2^\mu \tau)$  sums to  $\mathcal{O}(\tau + \log(\ell/\tau))$ .

Evaluating  $\text{check}(i, j, 2^\mu \tau)$  requires computing  $\phi(T[i \dots i + 2^\mu \tau])$  and  $\phi(T[j \dots j + 2^\mu \tau])$ . The first fingerprint can be computed in constant time because  $i$  and  $i + 2^\mu \tau$  are always block aligned (see Lemma 2). The time to compute the second fingerprint depends on how far  $j$  and  $j + 2^\mu \tau$  each are from a sampled position, which in turn depends inversely on the significance of the block containing those positions. By Lemma 3,  $\mu$  is always a lower bound on the significance of  $j$ , which implies that  $\mu$  also lower bounds the significance of  $j + 2^\mu \tau$ , and thus by the way we sample positions, neither will have distance more than  $\tau/2^{\lfloor \mu/2 \rfloor}$  to a sampled position in  $\mathcal{S}$ . Finally, note that by the way  $\mu$  is increased and decreased,  $\text{check}(i, j, 2^\mu \tau)$  is called at most three times for any fixed value of  $\mu$ . Hence, the total time to compute all necessary fingerprints can be bounded as

$$\mathcal{O} \left( \sum_{\mu=0}^{\lg(\ell/\tau)} 1 + \tau/2^{\lfloor \mu/2 \rfloor} \right) = \mathcal{O}(\tau + \log(\ell/\tau)) .$$

### 3.5 The Las Vegas Data Structure

We now describe an  $\mathcal{O}(n^{3/2})$ -time and  $\mathcal{O}(n/\tau)$ -space algorithm for verifying that  $\phi$  is *collision-free* on all pairs of substrings of  $T$  that the query algorithm compares. If a collision is found we pick a new  $\phi$  and try again. With high probability we can find a collision-free  $\phi$  in a constant number of trials, so we obtain the claimed Las Vegas data structure.

If  $\tau \leq \sqrt{n}$  we use the verification algorithm of Bille et al. [2], using  $\mathcal{O}(n\tau + n \log n)$  time and  $\mathcal{O}(n/\tau)$  space. Otherwise, we use the simple  $\mathcal{O}(n^2/\tau)$ -time and  $\mathcal{O}(n/\tau)$ -space algorithm described below.

Recall that all fingerprint comparisons in our algorithm are of the form

$$\phi(T[k\tau \dots k\tau + 2^l\tau - 1]) \stackrel{?}{=} \phi(T[j \dots j + 2^l\tau - 1])$$

for some  $k \in [n/\tau], j \in [n], l \in [\log(n/\tau)]$ . The algorithm checks each  $l \in [\log(n/\tau)]$  separately. For a fixed  $l$  it stores the fingerprints  $\phi(T[k\tau \dots k\tau + 2^l\tau])$  for all  $k \in [n/\tau]$  in a hash table  $\mathcal{H}$ . This can be done in  $\mathcal{O}(n)$  time and  $\mathcal{O}(n/\tau)$  space. For every  $j \in [n]$  the algorithm then checks whether  $\phi(T[j \dots j + 2^l\tau]) \in \mathcal{H}$ , and if so, it verifies that the underlying two substrings are in fact the same by comparing them character by character in  $\mathcal{O}(2^l\tau)$  time. By maintaining the fingerprint inside a sliding window of length  $2^l\tau$ , the verification time for a fixed  $l$  becomes  $\mathcal{O}(n2^l\tau)$ , i.e.,  $\mathcal{O}(n^2/\tau)$  time for all  $l \in [\log(n/\tau)]$ .

### 3.6 Queries with Long LCEs

In this section we describe an  $\mathcal{O}(\frac{n}{\tau})$  space data structure that in constant time either correctly computes  $\text{LCE}(i, j)$  or determines that  $\text{LCE}(i, j) \leq \tau^2$ . The data structure can be constructed in  $\mathcal{O}(n \log \frac{n}{\tau})$  time by a Monte Carlo or Las Vegas algorithm.

**The Data Structure** Let  $\mathcal{S}_\tau \subseteq [1, n]$  called the *sampled positions* of  $T$  (to be defined below), and consider the sets  $A$  and  $B$  of suffixes of  $T$  and  $T^R$ , respectively.

$$A = \{T[i \dots] \mid i \in \mathcal{S}_\tau\} \quad , \quad B = \{T[\dots i]^R \mid i \in \mathcal{S}_\tau\} .$$

We store a data structure for  $A$  and  $B$ , that allows us to perform constant time longest common extension queries on any pair of suffixes in  $A$  or any pair in  $B$ . This can be achieved by well-known techniques, e.g., storing a sparse suffix tree for  $A$  and  $B$ , equipped with a nearest common ancestor data structure. To define  $\mathcal{S}_\tau$ , let  $D_\tau = \{0, 1, \dots, \tau\} \cup \{2\tau, \dots, (\tau-1)\tau\}$ , then

$$\mathcal{S}_\tau = \{1 \leq i \leq n \mid i \bmod \tau^2 \in D_\tau\} . \tag{1}$$

**Answering a Query** To answer an LCE query, we need the following definitions. For  $i, j \in \mathcal{S}_\tau$  let  $\text{LCE}_R(i, j)$  denote the longest common prefix of  $T[\dots i]^R \in B$  and  $T[\dots j]^R \in B$ . Moreover, for  $i, j \in [n]$ , we define the function

$$\delta(i, j) = (((i - j) \bmod \tau) - i) \bmod \tau^2. \quad (2)$$

We will write  $\delta$  instead of  $\delta(i, j)$  when  $i$  and  $j$  are clear from the context.

The following lemma gives the key property that allows us to answer a query.

**Lemma 4.** *For any  $i, j \in [n - \tau^2]$ , it holds that  $i + \delta, j + \delta \in \mathcal{S}_\tau$ .*

*Proof.* Direct calculation shows that  $(i + \delta) \bmod \tau^2 \leq \tau$ , and that  $(j + \delta) \bmod \tau = 0$ , and thus by definition both  $i + \delta$  and  $j + \delta$  are in  $\mathcal{S}_\tau$ .

To answer a query  $\text{LCE}(i, j)$ , we first verify that  $i, j \in [n - \tau^2]$  and that  $\text{LCE}_R(i + \delta, j + \delta) \geq \delta$ . If this is not the case, we have established that  $\text{LCE}(i, j) \leq \delta < \tau^2$ , and we stop. Otherwise, we return  $\delta + \text{LCE}(i + \delta, j + \delta) - 1$ .

**Analysis** To prove the correctness, suppose  $i, j \in [n - \tau^2]$  (if not clearly  $\text{LCE}(i, j) < \tau^2$ ) then we have that  $i + \delta, j + \delta \in \mathcal{S}_\tau$  (Lemma 4). If  $\text{LCE}_R(i + \delta, j + \delta) \geq \delta$  it holds that  $T[i \dots i + \delta] = T[j \dots j + \delta]$  so the algorithm correctly computes  $\text{LCE}(i, j)$  as  $\delta + 1 + \text{LCE}(i + \delta, j + \delta)$ . Conversely, if  $\text{LCE}_R(i + \delta, j + \delta) < \delta$ ,  $T[i \dots i + \delta] \neq T[j \dots j + \delta]$  it follows that  $\text{LCE}(i, j) < \delta < \tau^2$ .

Query time is  $\mathcal{O}(1)$ , since computing  $\delta$ ,  $\text{LCE}_R(i + \delta, j + \delta)$  and  $\text{LCE}(i + \delta, j + \delta)$  all takes constant time. Storing the data structures for  $A$  and  $B$  takes space  $\mathcal{O}(|A| + |B|) = \mathcal{O}(|\mathcal{S}_\tau|) = \mathcal{O}(\frac{n}{\tau})$ . For the preprocessing stage, we can use recent algorithms by I et al. [8] for constructing the sparse suffix tree for  $A$  and  $B$  in  $\mathcal{O}(\frac{n}{\tau})$  space. They provide a Monte Carlo algorithm using  $\mathcal{O}(n \log \frac{n}{\tau})$  time (correct w.h.p.), and a Las Vegas algorithm using  $\mathcal{O}(\frac{n}{\tau})$  time (w.h.p.).

## 4 Derandomizing the Monte Carlo Data Structure

Here we give a general technique for derandomizing Karp-Rabin fingerprints, and apply it to our Monte Carlo algorithm. The main result is that for any constant  $\varepsilon > 0$ , the data structure can be constructed completely deterministically in  $\mathcal{O}(n^{2+\varepsilon})$  time using  $\mathcal{O}(n/\tau)$  space. Thus, compared to the probabilistic preprocessing of the Las Vegas structure using  $\mathcal{O}(n^{3/2})$  time with high probability, it is relatively cheap to derandomize the data structure completely.

Our derandomizing technique is stated in the following lemma.

**Lemma 5.** *Let  $A, L \subset \{1, 2, \dots, n\}$  be a set of positions and lengths respectively such that  $\max(L) = n^{\Omega(1)}$ . For every  $\varepsilon \in (0, 1)$ , there exist a fingerprinting function  $\phi$  that can be evaluated in  $\mathcal{O}(\frac{1}{\varepsilon})$  time and has the property that for all  $a \in A, l \in L, i \in \{1, 2, \dots, n\}$ :*

$$\phi(T[a \dots a + (l - 1)]) = \phi(T[i \dots i + (l - 1)]) \iff T[a \dots a + (l - 1)] = T[i \dots i + (l - 1)]$$

We can find such a  $\phi$  using  $\mathcal{O}\left(\frac{S}{\varepsilon}\right)$  space and  $\mathcal{O}\left(\frac{n^{1+\varepsilon} \log n}{\varepsilon^2} \frac{|A|}{S} \max(L) |L|\right)$  time, for any value of  $S \in [1, |A|]$ .

The proof of the lemma appears in the full version of this paper [3].

**Corollary 2.** For any  $\tau \in [1, n]$ , the LCE problem can be solved by a deterministic data structure with  $\mathcal{O}(n/\tau)$  space usage and  $\mathcal{O}(\tau)$  query time. The data structure can be constructed in  $\mathcal{O}(n^{2+\varepsilon})$  time using  $\mathcal{O}(n/\tau)$  space.

*Proof.* We use the lemma with  $A = \{k\tau \mid k \in [n/\tau]\}$ ,  $L = \{2^l \tau \mid l \in [\log(n/\tau)]\}$ ,  $S = |A| = n/\tau$  and a suitable small constant  $\varepsilon > 0$ .

## References

1. A. Amir, M. Lewenstein, and E. Porat. Faster algorithms for string matching with  $k$  mismatches. *J. Algorithms*, 50(2):257–275, 2004.
2. P. Bille, I. L. Gørtz, B. Sach, and H. W. Vildhøj. Time-space trade-offs for longest common extensions. *J. of Discrete Algorithms*, 25:42–50, 2014.
3. P. Bille, I. L. Gørtz, M. B. T. Knudsen, M. Lewenstein, and H. W. Vildhøj. Longest Common Extensions in Sublinear Space. *arXiv:1504.02671*, 2015.
4. R. Cole and R. Hariharan. Approximate String Matching: A Simpler Faster Algorithm. *SIAM J. Comput.*, 31(6):1761–1782, 2002.
5. N. J. Fine and H. S. Wilf. Uniqueness Theorems for Periodic Functions. *Proc. AMS*, 16(1):109–114, 1965.
6. D. Gusfield and J. Stoye. Linear time algorithms for finding and representing all the tandem repeats in a string. *J. Comput. Syst. Sci.*, 69:525–546, 2004.
7. D. Harel and R. E. Tarjan. Fast algorithms for finding nearest common ancestors. *SIAM J. Comput.*, 13(2):338–355, 1984.
8. T. I, J. Kärkkäinen, and D. Kempa. Faster Sparse Suffix Sorting. In *Proc. 31st STACS*, volume 25, pages 386–396, Dagstuhl, Germany, 2014.
9. R. M. Karp and M. O. Rabin. Efficient randomized pattern-matching algorithms. *IBM J. Res. Dev.*, 31(2):249–260, 1987.
10. R. Kolpakov and G. Kucherov. Searching for gapped palindromes. In *Proc. 19th CPM, LNCS*, volume 5029, pages 18–30, 2008.
11. G. M. Landau, E. W. Myers, and J. P. Schmidt. Incremental string comparison. *SIAM J. Comput.*, 27(2):557–582, 1998.
12. G. M. Landau and J. P. Schmidt. An Algorithm for Approximate Tandem Repeats. *J. Comput. Biol.*, 8(1):1–18, 2001.
13. G. M. Landau and U. Vishkin. Fast Parallel and Serial Approximate String Matching. *J. Algorithms*, 10:157–169, 1989.
14. M. G. Main and R. J. Lorentz. An  $O(n \log n)$  algorithm for finding all repetitions in a string. *J. Algorithms*, 5(3):422–432, 1984.
15. G. Manacher. A New Linear-Time “On-Line” Algorithm for Finding the Smallest Initial Palindrome of a String. *J. ACM*, 22(3):346–351, 1975.
16. E. W. Myers. An  $O(ND)$  difference algorithm and its variations. *Algorithmica*, 1(2):251–266, 1986.
17. P. Weiner. Linear pattern matching algorithms. In *Proc. 14th FOCS (SWAT)*, pages 1–11, 1973.